# Building Robust Datasets for Commonsense Reasoning

Rowan Zellers

rowanzellers.com

A man has a few tools and is pumping his car up so he can take off the tire.
He...

A man has a few tools and is pumping his car up so he can take off the tire.
He...

A. stops on the front of the bike and moves it to the left.

B. gets out of the car while leaving the engine running.

C. uses the tool to take off all of the nuts one by one.

D. goes down from the cars, landing straight in.

A man has a few tools and is pumping his car up so he can take off the tire.
He...

A. stops on the front of the bike and moves it to the left.

B. gets out of the car while leaving the engine running.

C. uses the tool to take off all of the nuts one by one.

D. goes down from the cars, landing straight in.

A man has a few tools and is pumping his car up so he can take off the tire.
He...



A. stops on the front of the bike and

one by one.

D. goes down from the cars, landing straight in.

**This is natural language inference that requires commonsense reasoning!**

# Our contributions with *SWAG*

- SWAG: Natural Language Inference + Commonsense Reasoning

# Our contributions with SWAG

- SWAG: Natural Language Inference + Commonsense Reasoning
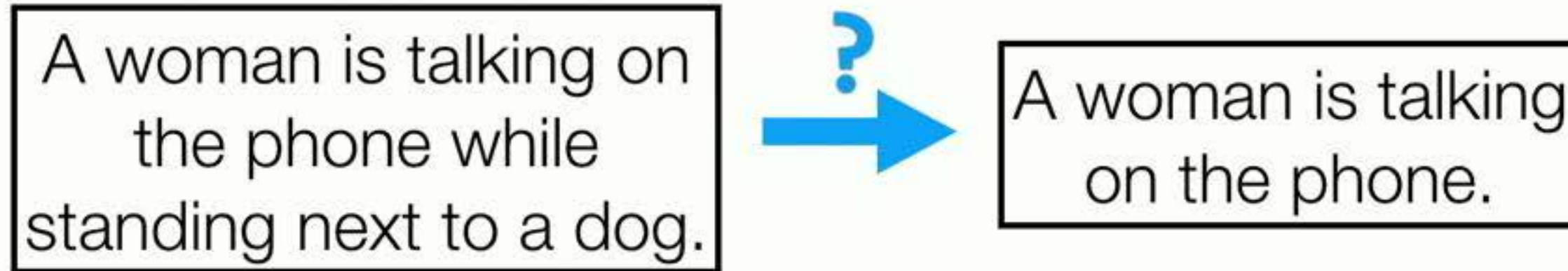
- Adversarial Filtering

# Scope of Natural Language Inference (NLI)

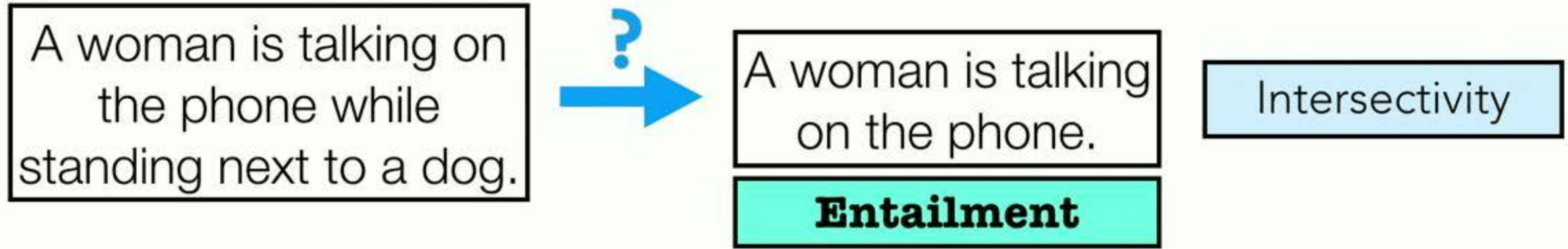Much of today's NLI requires only linguistic knowledge, without as much commonsense reasoning.

*Bowman et al.,
2015,
Wang et al., 2018*

# Scope of Natural Language Inference (NLI)

Much of today's NLI requires only linguistic knowledge, without as much commonsense reasoning. **Examples from SNLI:**
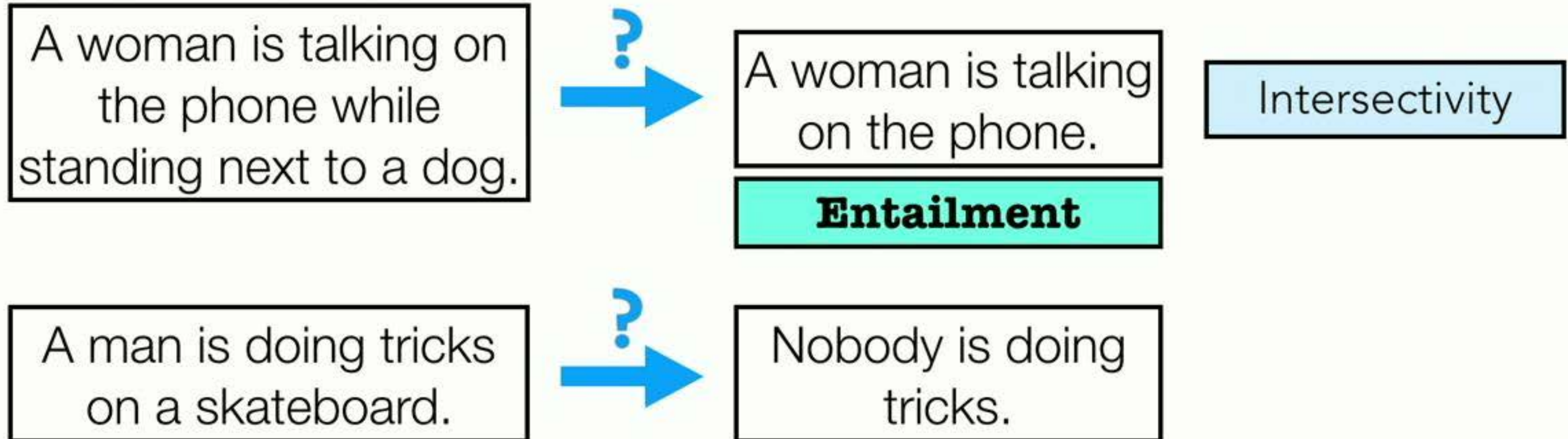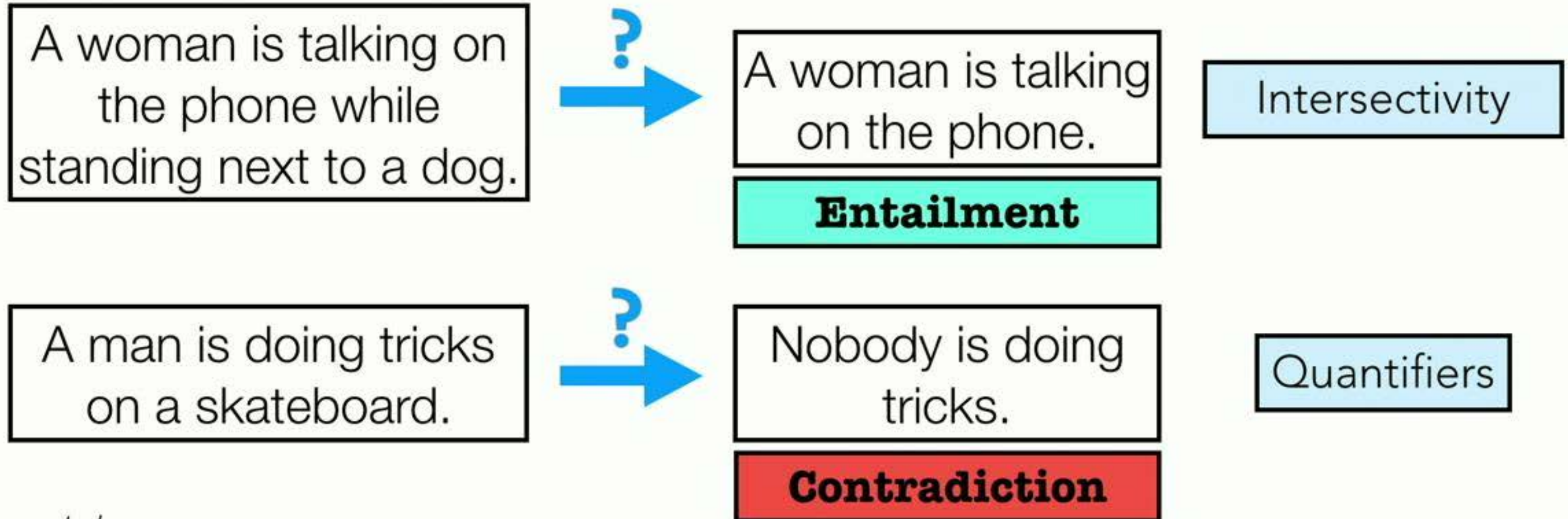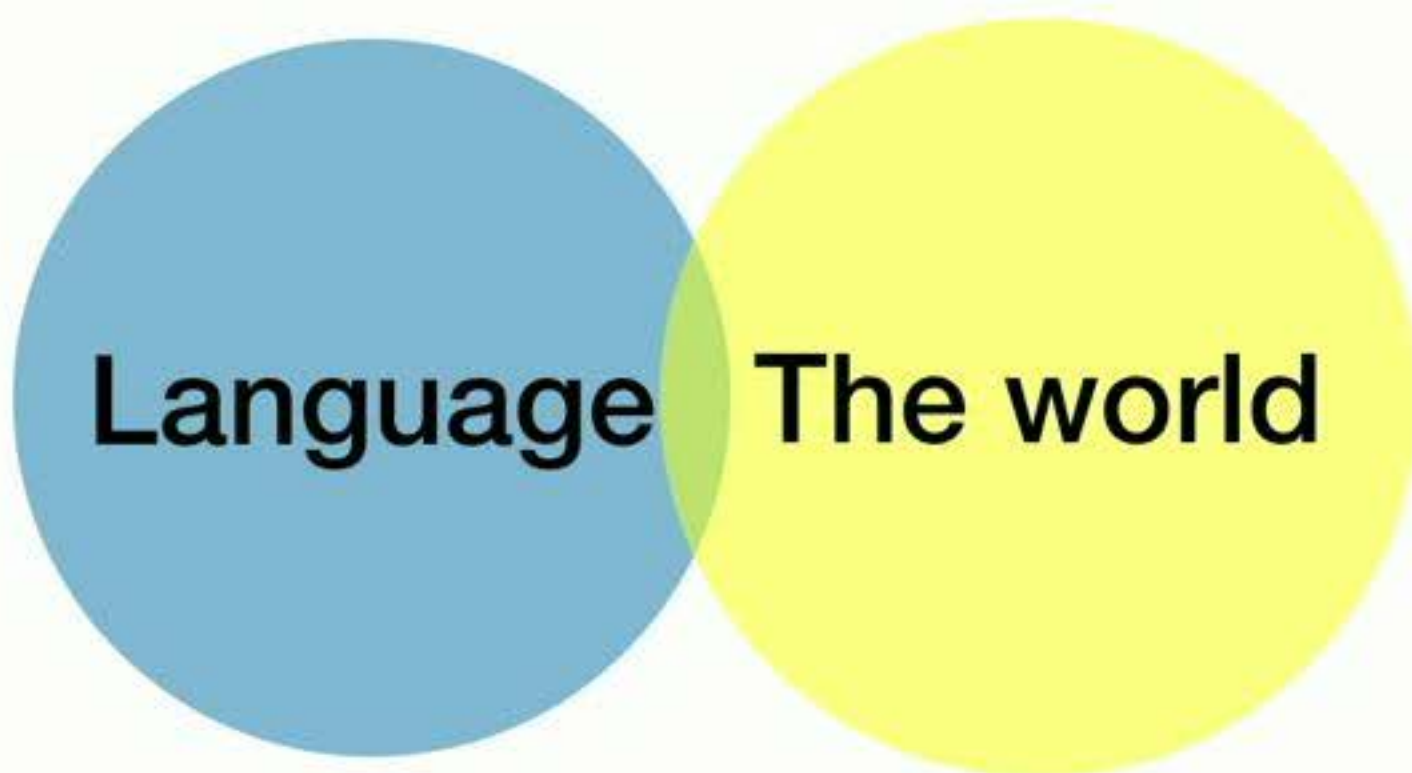
| A woman is talking on the phone while standing next to a dog. | **?** → | A woman is talking on the phone. |
|---|---|---|

*Bowman et al., 2015,
Wang et al., 2018*

# Scope of Natural Language Inference (NLI)

Much of today's NLI requires only linguistic knowledge, without as much commonsense reasoning. **Examples from SNLI:**

| A woman is talking on the phone while standing next to a dog. | ? → | A woman is talking on the phone. | Intersectivity |
| --- | --- | --- | --- |
| | | **Entailment** | |

Bowman et al., 2015,
Wang et al., 2018

# Scope of Natural Language Inference (NLI)

Much of today's NLI requires only linguistic knowledge, without as much commonsense reasoning. **Examples from SNLI:**

| A woman is talking on the phone while standing next to a dog. | ? → | A woman is talking on the phone. | Intersectivity |
| --- | --- | --- | --- |
| | | **Entailment** | |

| A man is doing tricks on a skateboard. | ? → | Nobody is doing tricks. |
| --- | --- | --- |

*Bowman et al., 2015,*
*Wang et al., 2018*

# Scope of Natural Language Inference (NLI)

Much of today's NLI requires only linguistic knowledge, without as much commonsense reasoning. **Examples from SNLI:**

| A woman is talking on the phone while standing next to a dog. | ? → | A woman is talking on the phone. | Intersectivity |
|---|---|---|---|
| | | **Entailment** | |
| A man is doing tricks on a skateboard. | ? → | Nobody is doing tricks. | Quantifiers |
| | | **Contradiction** | |

Bowman et al., 2015,
Wang et al., 2018

# Re-emphasizing Commonsense in NLI

Natural Language Inference
was supposed to be:

Language  The world

(Dagan et. al, 2006; LoBue and Yates, 2011)

# Situations with Adversarial Generations (*SWAG*)

A man has a few tools and is pumping his car up so he can take off the tire. He



t

*Krishna et al., 2017,
Rohrbach et al., 2016*

# Situations with Adversarial Generations (*SWAG*)



A man has a few tools and is pumping his car up so he can take off the tire. He

A. stops on the front of the bike and moves it to the left.

B. gets out of the car while leaving the engine running.

C. uses the tool to take off all of the nuts one by one.

D. goes down from the cars, landing straight in

ACTIVITYNET
LSMDC

*Krishna et al., 2017,*
*Rohrbach et al., 2016*

# Situations with Adversarial Generations (*SWAG*)

A man has a few tools and is pumping his car up so he can take off the tire. He

A. stops on the front of the bike and moves it to the left.

B. gets out of the car while leaving the engine running.

C. uses the tool to take off all of the nuts one by one.

D. goes down from the cars, landing straight in

**SWAG is purely a natural language inference task: the video is never used.**

t

t+1

Krishna et al., 2017, Rohrbach et al., 2016

ACTIVITYNET

LSMDC

# Situations with Adversarial Generations (*SWAG*)

A man has a few tools and is pumping his car up so he can take off the tire. He

A. stops on the front of the bike and moves it to the left.

B.

**How do we get the wrong answers?**

C.

nuts one by one.

D. goes down from the cars, landing straight in

ly a natural rence task: ever used.

# Our contributions

- SWAG: Natural Language Inference + Commonsense Reasoning

- Adversarial Filtering

# Annotation Artifacts

Human written datasets are susceptible to *annotation artifacts*: stylistic patterns that give unwanted clues for the labels.

*Schwartz et al., 2017*
*Gururangan et al., 2018*
*Poliak et al., 2018*

# Annotation Artifacts

Human written datasets are susceptible to *annotation artifacts*: stylistic patterns that give unwanted clues for the labels.

Schwartz et al., 2017
Gururangan et al., 2018
Poliak et al., 2018

A man is doing tricks on a skateboard.

**?**

Nobody is doing tricks.

**Contradiction**

# Annotation Artifacts

Human written datasets are susceptible to *annotation artifacts*: stylistic patterns that give unwanted clues for the labels.

*Schwartz et al., 2017*
*Gururangan et al., 2018*
*Poliak et al., 2018*

A man is doing tricks on a skateboard.

**?** →

Nobody is doing tricks.

**Contradiction**

For *Contradiction*, "Nobody" works...

# Annotation Artifacts

Human written datasets are susceptible to *annotation artifacts*: stylistic patterns that give unwanted clues for the labels.

*Schwartz et al., 2017*
*Gururangan et al., 2018*
*Poliak et al., 2018*

A man is doing tricks on a skateboard.

**Nobody** is doing tricks.

**Contradiction**

For *Contradiction*, "Nobody" works...

Nobody -> Negation

# A different approach...

# A different approach...

Let's have machines write the wrong endings!

# A different approach...

Let's have machines write the wrong endings!

Humans will **verify** that the original endings are better than the wrong ones.

# Adversarial Filtering



We'll *massively* *oversample* candidate endings.

# Adversarial Filtering



We'll *massively oversample* candidate endings.



And we'll use more models to remove obvious artifacts!

# Adversarial Filtering

data

We'll *massively oversample* candidate endings.

And we'll use more models to remove obvious artifacts!

# Adversarial Filtering

A man has a few tools and is pumping his car up so he can take off the tire.

He  uses the tool to take off all of the nuts one by one.

**NP**

**Ground truth VP**

# Adversarial Filtering

A man has a few tools and is pumping his car up so he can take off the tire.

He **uses the tool to take off all of the nuts one by one.**

*NP*

*Ground truth VP*

*LM generated endings*

- goes out onto the street.
- sits looking at the man.
- goes down from the cars, landing straight in.
- hauls the car over the bridge.

....

# Adversarial Filtering



**Ground truth VP** — uses the tool to take off all of the nuts one by one. → $x_1^+$

**LM generated endings**
- goes out onto the street.
- sits looking at the man.
- goes down from the cars, landing straight in.
- hauls the car over the bridge.
  ....

$x_{1,1}^-$
$x_{1,2}^-$
$x_{1,3}^-$
$x_{1,4}^-$

# Adversarial Filtering

# What AF does and doesn't do

# What AF does and doesn't do

Train $\sim$ Test ?

# What AF does and doesn't do

# What AF does and doesn't do

Train $\neq$ Test ?

# What AF does and doesn't do

THIS MODEL SHOULD ONLY
BE EVALUATED ON DATA
THAT'S SIMILAR TO THE
TRAINING SET!!!!

# What AF does and doesn't do

# What AF does and doesn't do

# Math behind AF

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{1 \leq i \leq N}$       Model $f_\theta : \mathcal{X} \to \mathbb{R}^{|\mathcal{Y}|}$

Loss function $L(f_\theta, \mathcal{D})$

# Math behind AF

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{1 \leq i \leq N}$  Model $f_\theta : \mathcal{X} \to \mathbb{R}^{|\mathcal{Y}|}$

Loss function $L(f_\theta, \mathcal{D})$

Empirical error

$$I(\mathcal{D}, f) = \frac{1}{N} \sum_{i=1}^{N} L(f_{\theta_i^\star}, \{(x_i, y_i)\}),$$

$$\text{where } \theta_i^\star = \operatorname*{argmin}_\theta L(f_\theta, \mathcal{D} \setminus \{(x_i, y_i)\}),$$

# Math behind AF

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{1 \leq i \leq N}$     Model $f_\theta : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$

Loss function $L(f_\theta, \mathcal{D})$

Empirical error

$$I(\mathcal{D}, f) = \frac{1}{N} \sum_{i=1}^{N} L(f_{\theta_i^\star}, \{(x_i, y_i)\}),$$

$$\text{where } \theta_i^\star = \underset{\theta}{\mathrm{argmin}} \, L(f_\theta, \mathcal{D} \setminus \{(x_i, y_i)\}),$$

The best parameters as determined by other datapoints

# Math behind AF

Dataset $\mathcal{D} = \{(x_i, y_i)\}_{1 \leq i \leq N}$     Model $f_\theta : \mathcal{X} \to \mathbb{R}^{|\mathcal{Y}|}$

How much does knowing N-1 datapoints help?

Empirical error

$$I(\mathcal{D}, f) = \frac{1}{N} \sum_{i=1}^{N} L(f_{\theta_i^\star}, \{(x_i, y_i)\}),$$

$$\text{where } \theta_i^\star = \underset{\theta}{\operatorname{argmin}} \, L(f_\theta, \mathcal{D} \setminus \{(x_i, y_i)\}),$$

The best parameters as determined by other datapoints

# Math behind AF

Let's **MAXIMIZE** dataset-level empirical error!

How much does knowing N-1 datapoints help?

Empirical error

$$I(\mathcal{D}, f) = \frac{1}{N} \sum_{i=1}^{N} L(f_{\theta_i^\star}, \{(x_i, y_i)\}),$$

$$\text{where } \theta_i^\star = \underset{\theta}{\operatorname{argmin}} \, L(f_\theta, \mathcal{D} \setminus \{(x_i, y_i)\}),$$

The best parameters as determined by other datapoints

# Adversarial Filters for SWAG



We used an ensemble of stylistic models centered around the ending sentence. These are powerful adversaries!

(Schwartz et al., 2017, Gururangan et al., 2018)

# Adversarial Filters for SWAG

We used an ensemble of stylistic models centered around the ending sentence. These are powerful adversaries!
(Schwartz et al., 2017, Gururangan et al., 2018)

A multilayer perceptron, given LM perplexities as features

Bag-of-words

CNN

BiLSTM, with uncommon words replaced by POS tags

# Adversarial Filters for SWAG

We used an ensemble of stylistic models centered around the ending sentence. These are powerful adversaries!
(Schwartz et al., 2017, Gururangan et al., 2018)

**After 100 iterations of adversarial filtering, the accuracy of these stylistic models drops to chance.**

# Adversarial Filters for SWAG

data

AF

We used an ensemble of stylistic models centered around the ending sentence. These are powerful adversaries!
(Schwartz et al., 2017, Gururangan et al., 2018)

**After 100 iterations of adversarial filtering, the accuracy of these stylistic models drops to chance.**

**Last, crowd workers validate the entire dataset to remove false negatives, leaving us with 110k multiple-choice questions.**

# Unique Contributions of Adversarial Filtering

# Unique Contributions of Adversarial Filtering



Diversity of generated+filtered sentences is limited only by your corpus

# Unique Contributions of Adversarial Filtering



Diversity of generated+filtered sentences is limited only by your corpus



Can raise the bar of difficulty during dataset construction

# Unique Contributions of Adversarial Filtering

Diversity of generated+filtered sentences is limited only by your corpus

Can raise the bar of difficulty during dataset construction

Human validation is cheap as workers don't write the endings

# SWAG Results

# SWAG Results

# SWAG Results

SWAG Results

# SWAG Results

# SWAG Results



A bar chart titled "SWAG Results" with y-axis labeled "Accuracy" ranging from 0 to 100. Bars: Human = 88, FastText = 30, LSTM+GloVe = 46, ESIM = 53, ESIM+ELMo = 59, OpenAI GPT = 78, BERT = (no value shown).

# SWAG Results

# What's next?

# What do models learn by fine-tuning on SWAG?



A girl is going across the monkey bars.

She gets to the other side and stands on a wooden plank.

# What do models learn by fine-tuning on SWAG?



A girl is going across the monkey bars.

She gets to the other side and stands on a wooden plank.

Can a 2 sentence NLI dataset be constructed at scale while being resistant to pretraining approaches?

Can a 2 sentence NLI dataset be constructed at scale while being resistant to pretraining approaches?

**Actually, yes.**

Someone is pointing at someone else. The waiter...

Why is [person4 🖼] pointing at [person1 🖼]?

How can we get vision systems to learn this??

Why is [person4 🖼] pointing at [person1 🖼]?

Why is [person4 🖼] pointing at [person1 🖼]?

a) **He is telling [person3 🖼] that [person1 🖼] ordered the pancakes.**

b) He just told a joke.

c) He is feeling accusatory towards [person1 🖼].

d) He is giving [person1 🖼] directions.

Many AI systems perform well, but do so for *questionable reasons*

Why is [person4 🧑] pointing at [person1 👤]?

**a) He is telling [person3 🧑] that [person1 🧑] ordered the pancakes.**

b) He just told a joke.

c) He is feeling accusatory towards [person1 🧑].

d) He is giving [person1 🧑] directions.

*I chose a) because…*

a) [person1 🧑] has the pancakes in front of him.

b) [person4 🧑] is taking everyone's order and asked for clarification.

c) [person3 🧑] is looking at the pancakes and both she and [person2 🧑] are smiling slightly.

**d) [person3 🧑] is delivering food to the table, and she might not know whose order is whose.**

# Our contributions



- New task: Visual Commonsense Reasoning

# Our contributions

- New task: Visual Commonsense Reasoning

- Building VCR, feat. Adversarial Matching

- Recognition to Cognition Networks

# Our contributions

- New task: Visual Commonsense Reasoning

- Building VCR, feat. Adversarial Matching

- Recognition to Cognition Networks

# Collecting commonsense inferences

# Collecting commonsense inferences



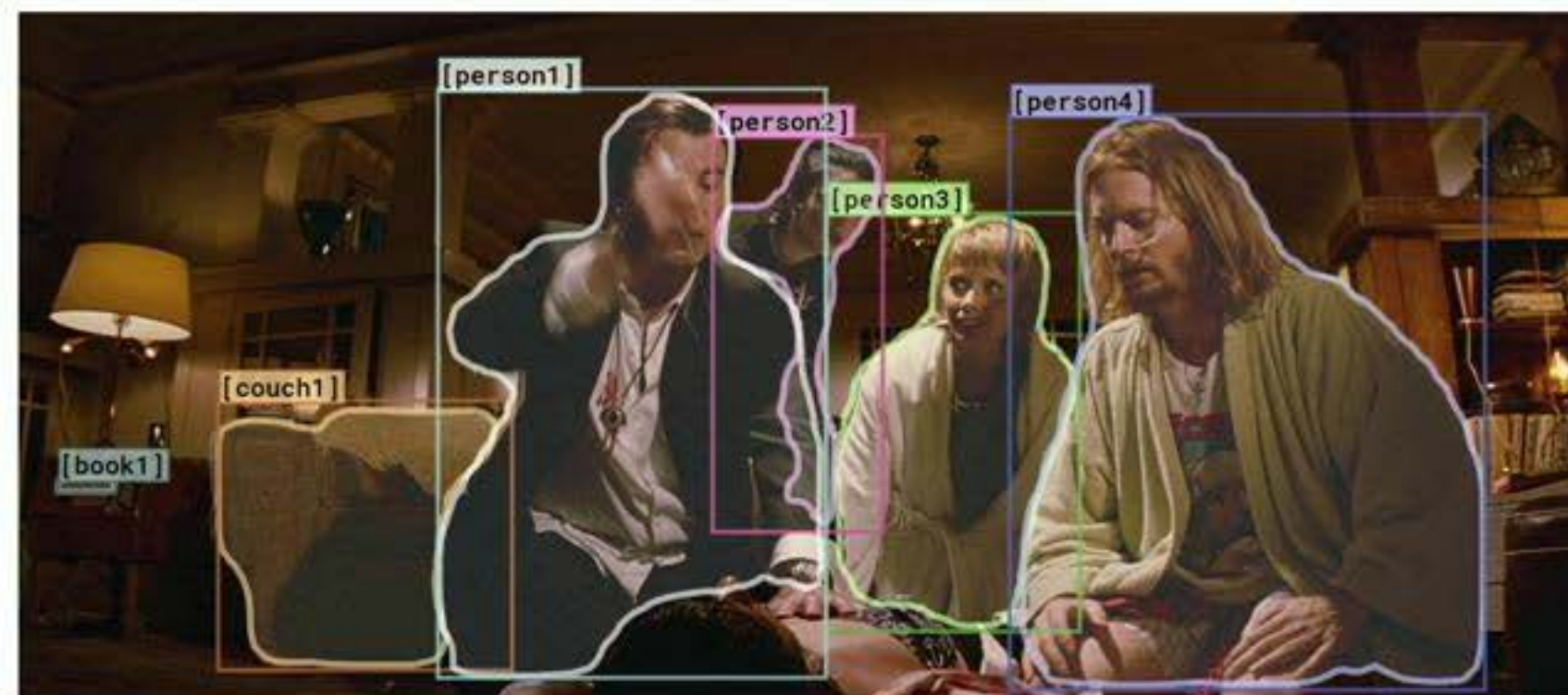fouhey et al, "from lifestyle vlogs to everyday interactions." cvpr2018

# Collecting commonsense inferences
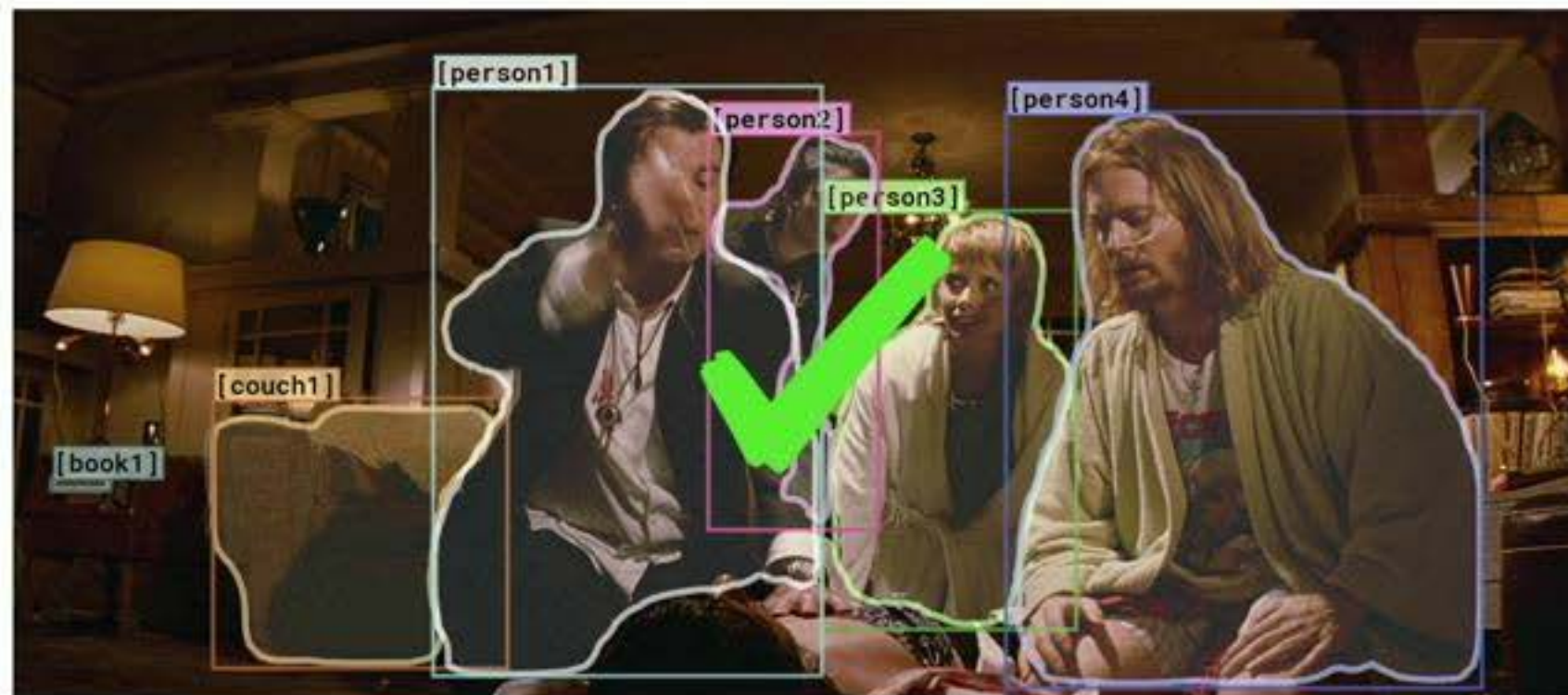
# Collecting commonsense inferences

# Collecting commonsense inferences

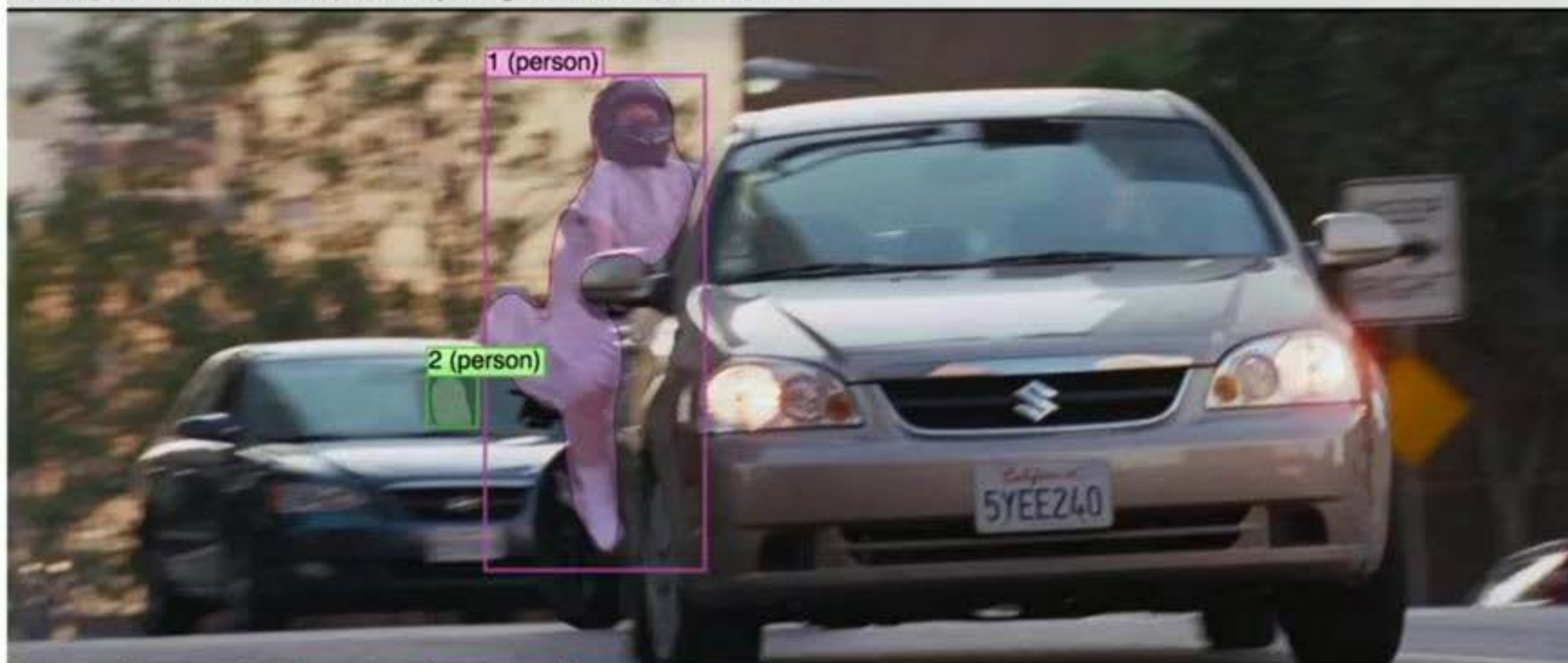# Collecting commonsense inferences

# Collecting commonsense inferences

# Collecting commonsense inferences



Examples (see rowanzellers.com/halloffame for more)   (expand/collapse)

Past caption: The driver looks out and sees the Ducati's front wheel waggling about in midair.

This caption: The bike's front wheel drops down, and SOMEONE speeds through town.

Next caption: SOMEONE weaves precariously through traffic and hops onto a pavement.

# Collecting commonsense inferences

# Collecting commonsense inferences

How do we get the wrong answers, avoiding annotation artifacts?

what color is the shirt?

what color is the shirt?

teal

blue

violet

magenta

# Adversarial Matching

Wrong answers must have be

relevant to the question   yet   different from the correct answer.

# Adversarial Matching

Wrong answers must have be

relevant to the question    yet    different from the correct answer.

Q,A'    **Question Relevance**

A,A'    **Entailment**

# Adversarial Matching

Wrong answers must have be

relevant to the question    yet    different from the correct answer.

*Q,A'*    **Question Relevance**

*A,A'*    **Entailment**

We'll use these two metrics to **recycle right answers** to other questions, using a minimum weight bipartite matching.
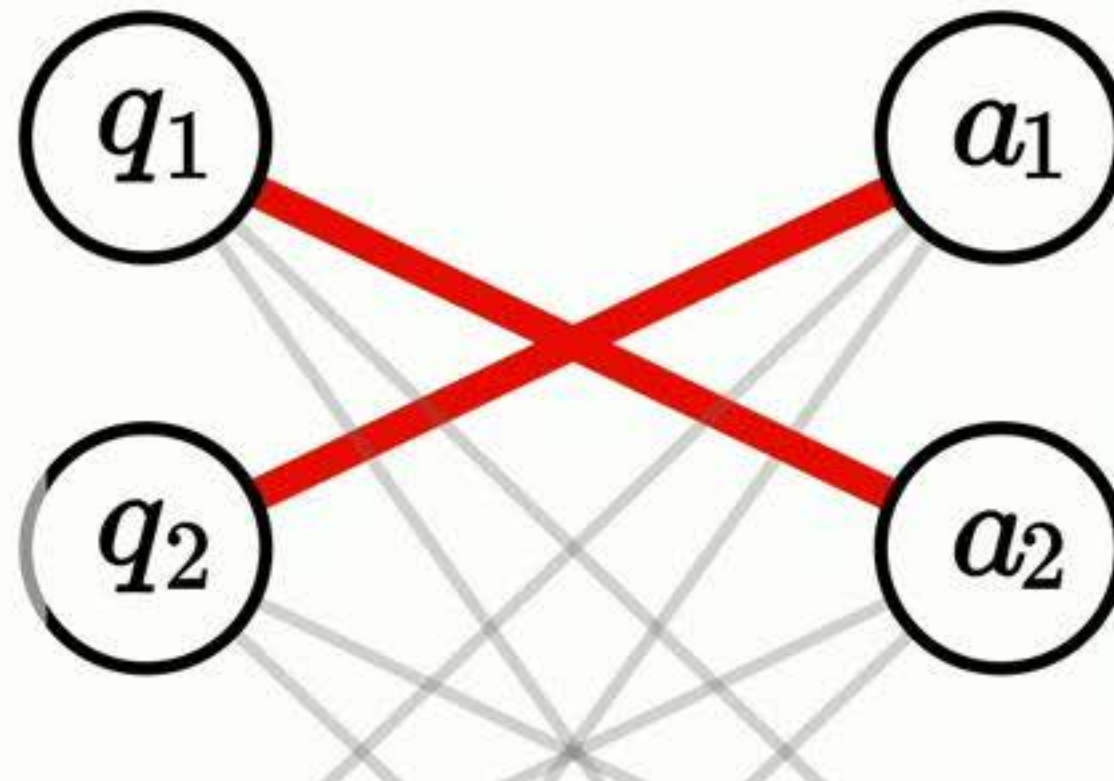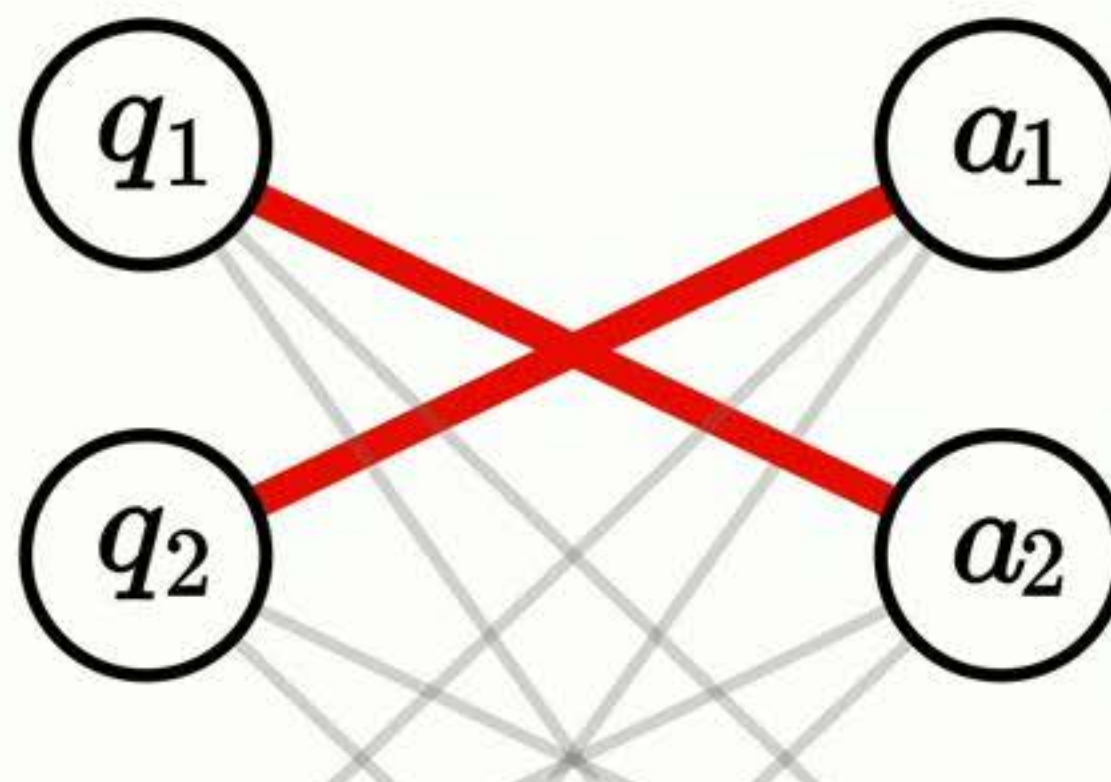
# Adversarial Matching

We'll use these two metrics to **recycle right answers** to other questions, using a minimum weight bipartite matching.

Why are [person1] and [person3] holding their foreheads together?

Why do [person1] and [person3] have their hands clasped?

$q_1$

$q_2$

$a_1$

$a_2$

They are about to kiss.

[person1] and [person3] are praying.

# Adversarial Matching

We'll use these two metrics to **recycle right answers** to other questions, using a minimum weight bipartite matching.

Why are [person1] and [person3] holding their foreheads together?

Why do [person1] and [person3] have their hands clasped?

$q_1$

$q_2$

$a_1$

$a_2$

They are about to kiss.

[person1] and [person3] are praying.

# Adversarial Matching

Wrong answers must have be

relevant to the question    yet    different from the correct answer.

Q,A'    **Question Relevance**

A,A'    **Entailment**

# What about the people tags (`[person5]`)?

# Adversarial Matching

Wrong answers must have be

relevant to the question    yet    different from the correct answer.

$Q,A'$ **Question Relevance**

$A,A'$ **Entailment**

Question relevance dominates - hard for machines

— Low $\lambda$

High $\lambda$—

Entailment penalty dominates - easy for humans

# Adversarial Matching

Wrong answers must have be

relevant to the question       yet       different from the correct answer.

Q,A'       *Question Relevance*                    A,A'       *Entailment*

This works for the rationales too!

Question relevance
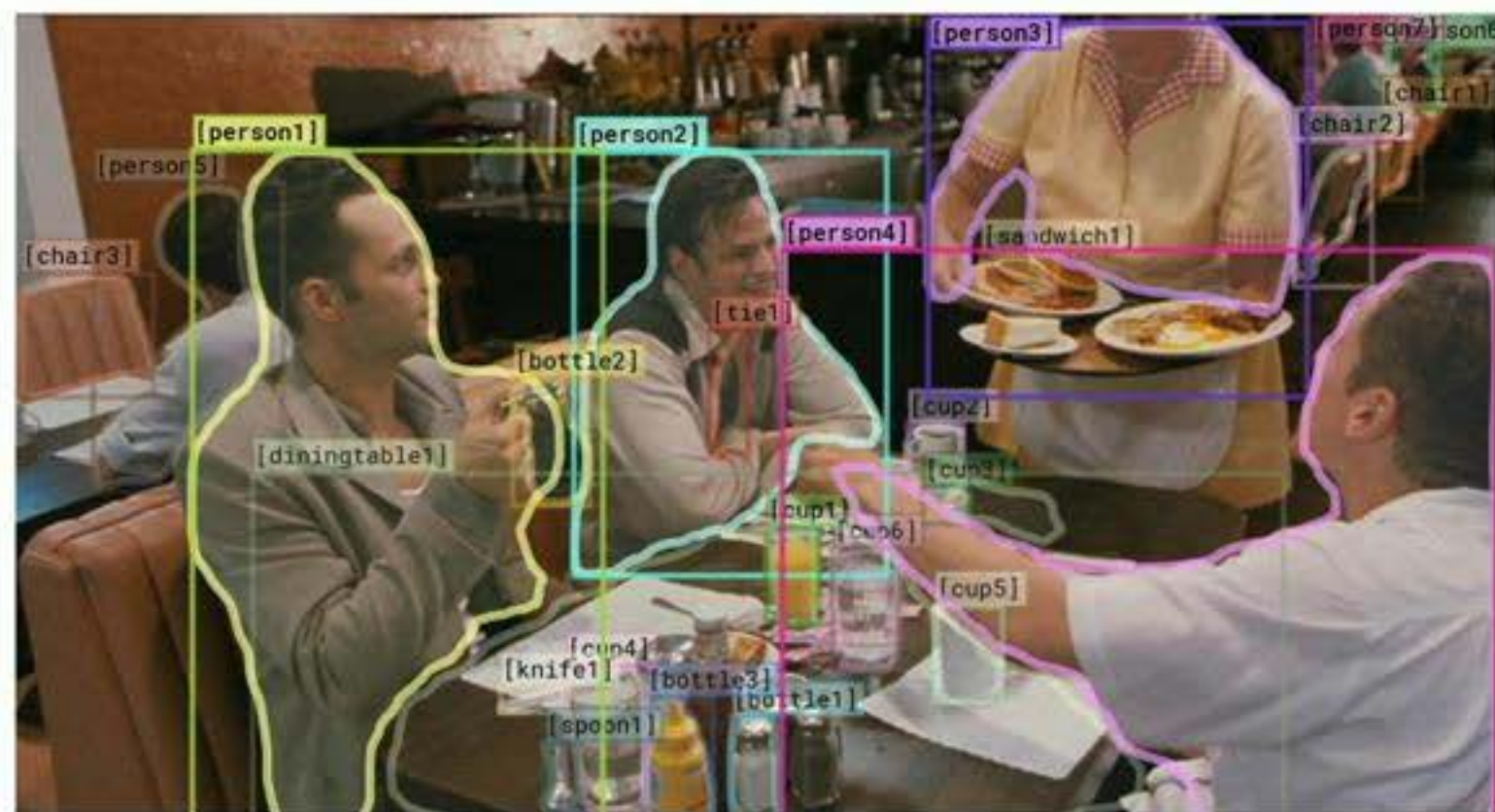dominates - hard
for machines

—Low λ

High λ—

Entailment
penalty dominates
- easy for humans

# What about the people tags (`[person5]`)?

We'll randomly modify the detection tags in candidate answer to better match the new question/image.
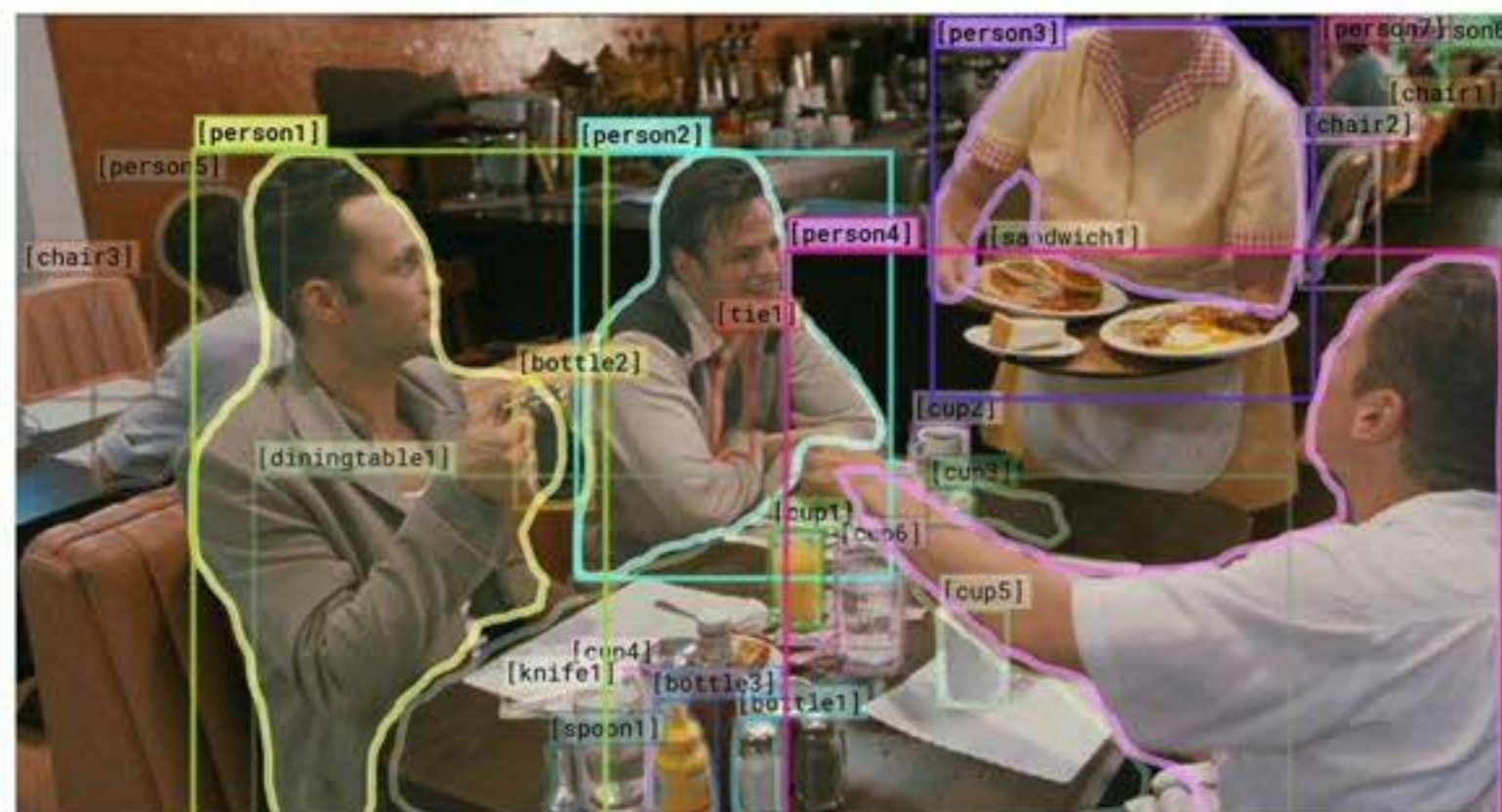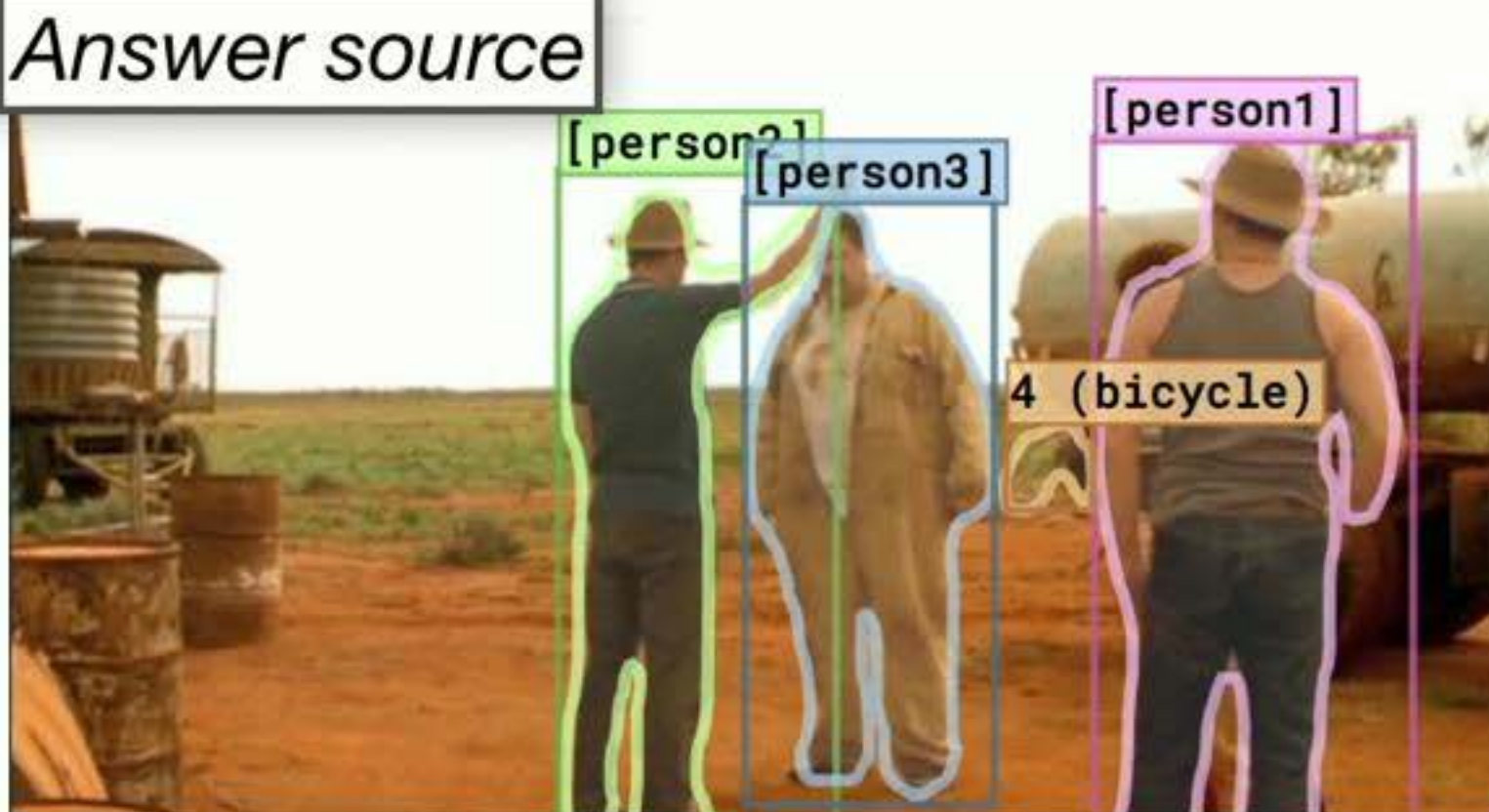
# What about the people tags (`[person5]`)?

We'll randomly modify the detection tags in candidate answer to better match the new question/image.

# What about the people tags ([person5])?

We'll randomly modify the detection tags in candidate answer to better match the new question/image.
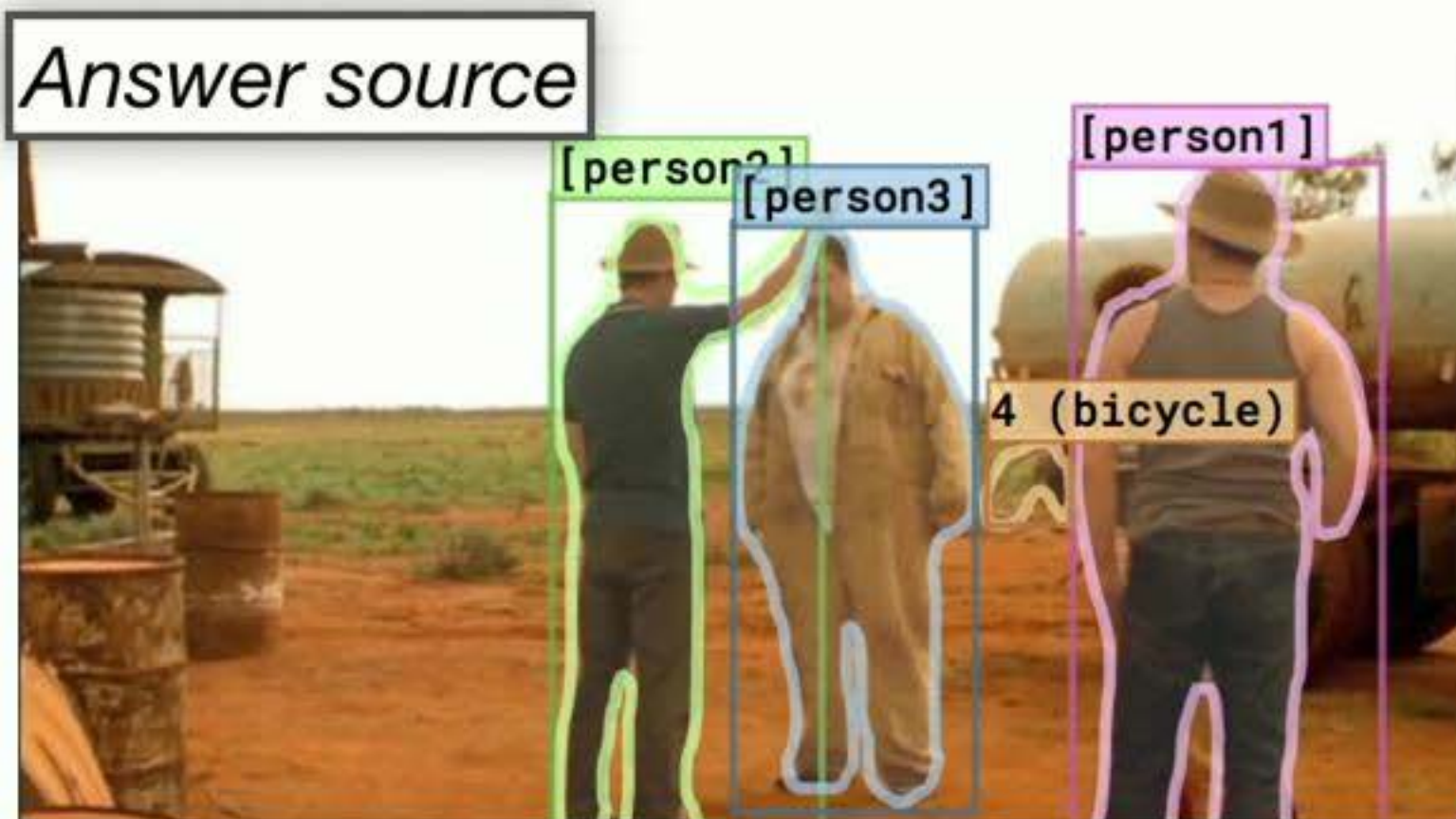


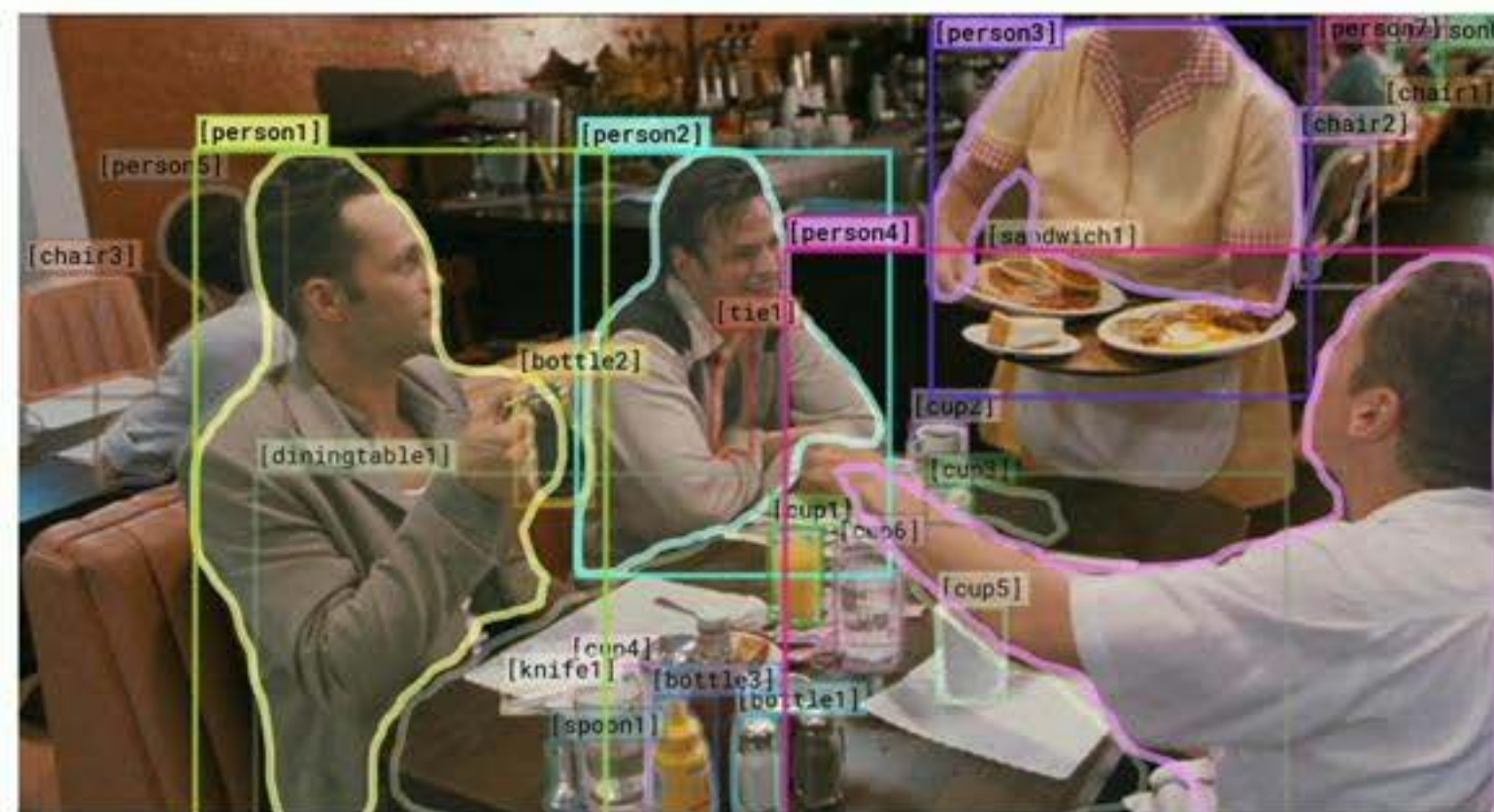Answer source

He is giving [person3 ] directions.

# What about the people tags ([person5])?

We'll randomly modify the detection tags in candidate answer to better match the new question/image.



Answer source

He is giving [person3 ] directions. ➡ He is giving [person1 ] directions.

# Unique Contributions of Adversarial Matching

No answer-only bias

**0.25**

# Unique Contributions of Adversarial Matching



No answer-only bias

0.25



Can raise the bar of difficulty during dataset construction

# Unique Contributions of Adversarial Matching



No answer-only bias

0.25



Can raise the bar of difficulty during dataset construction



Entailment NLI means that human validation isn't (as) needed

# Putting it all together: VCR

- VCR features 290k questions over 110k images, each with answers and rationales.

- The questions are diverse and challenging.

# Our contributions



- New task: Visual Commonsense Reasoning

- Building VCR, feat. Adversarial Matching

- Recognition to Cognition Networks

# Detour: Setting up the task

# Detour: Setting up the task

*Question answering*          *Answer justification*

## Question answering

$$Q \rightarrow A$$

## Answer justification

$$QA \rightarrow R$$

# Detour: Setting up the task

### Question answering

$$Q \rightarrow A$$

Query

### Answer justification

$$QA \rightarrow R$$

Query

**Let's make both tasks have the same format!**

# Detour: Setting up the task

*Question answering*

$$Q \rightarrow A$$

**Query**  **Responses**

*Answer justification*

$$QA \rightarrow R$$

**Query**  **Responses**

Let's make both tasks have the same format!

# Detour: Setting up the task

*Question answering*

$$Q \rightarrow A$$

Query     Responses

*Answer justification*

$$QA \rightarrow R$$

Query     Responses

$$Q \rightarrow AR$$

# Recognition to Cognition Networks

# Recognition to Cognition Networks

# Recognition to Cognition Networks

1. Grounding
2. Contextualization
3. Reasoning

*Part 1: Grounding*

Objects

*Query*

Why is [person4 🖼️] pointing at [person1 🖼️]?

*Response Choice*

He is telling [person3 🖼️] that [person1 🖼️] ordered pancakes.

*Part 1: Grounding*

Objects

Query

Why is [person4 ] pointing at [person1 ]?

Response Choice

He is telling [person3 ] that [person1 ] ordered pancakes.

Objects

*Part 1: Grounding*

*Query*

Why is [person4] pointing at [person1]?

*Response Choice*

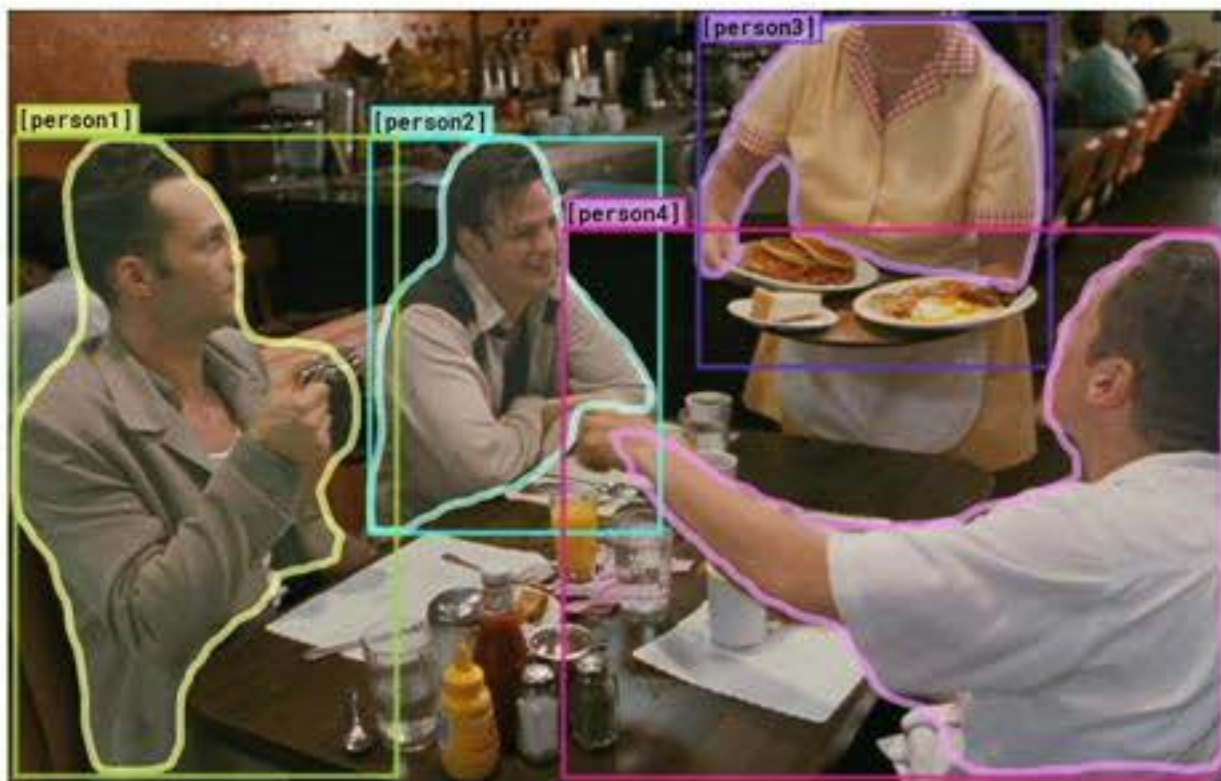He is telling [person3] that [person1] ordered pancakes.
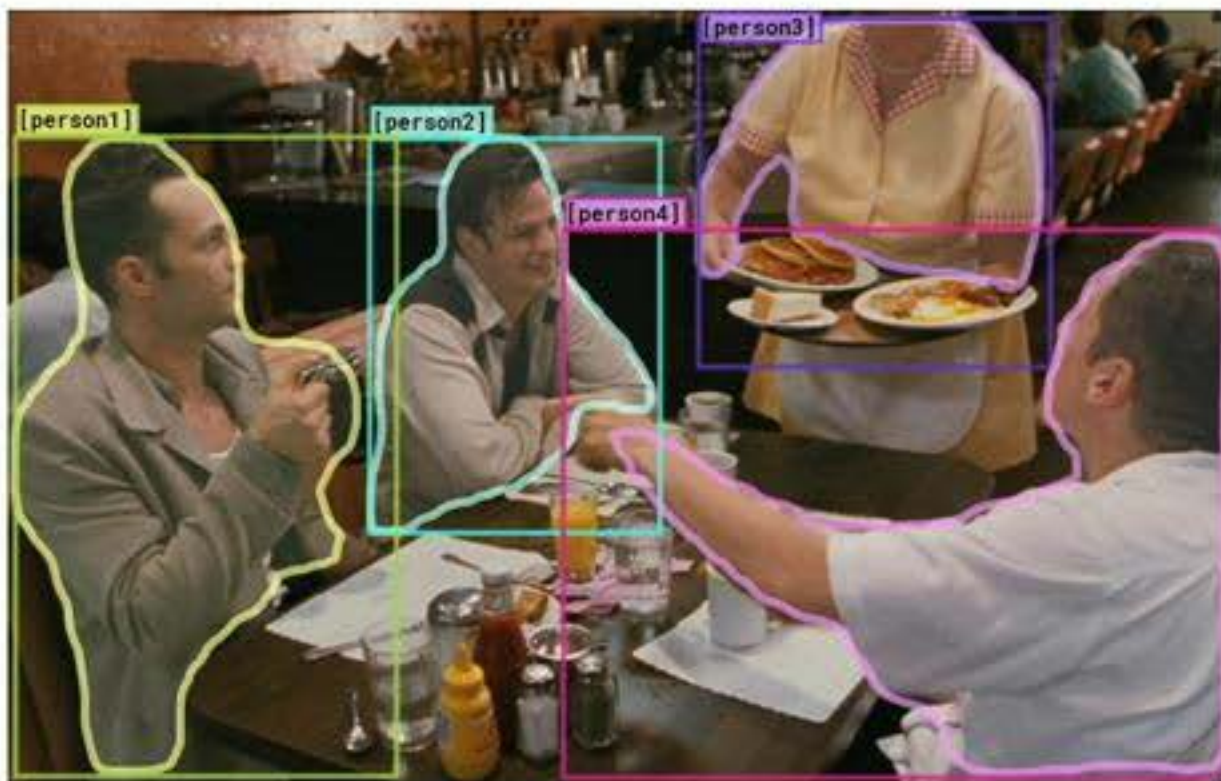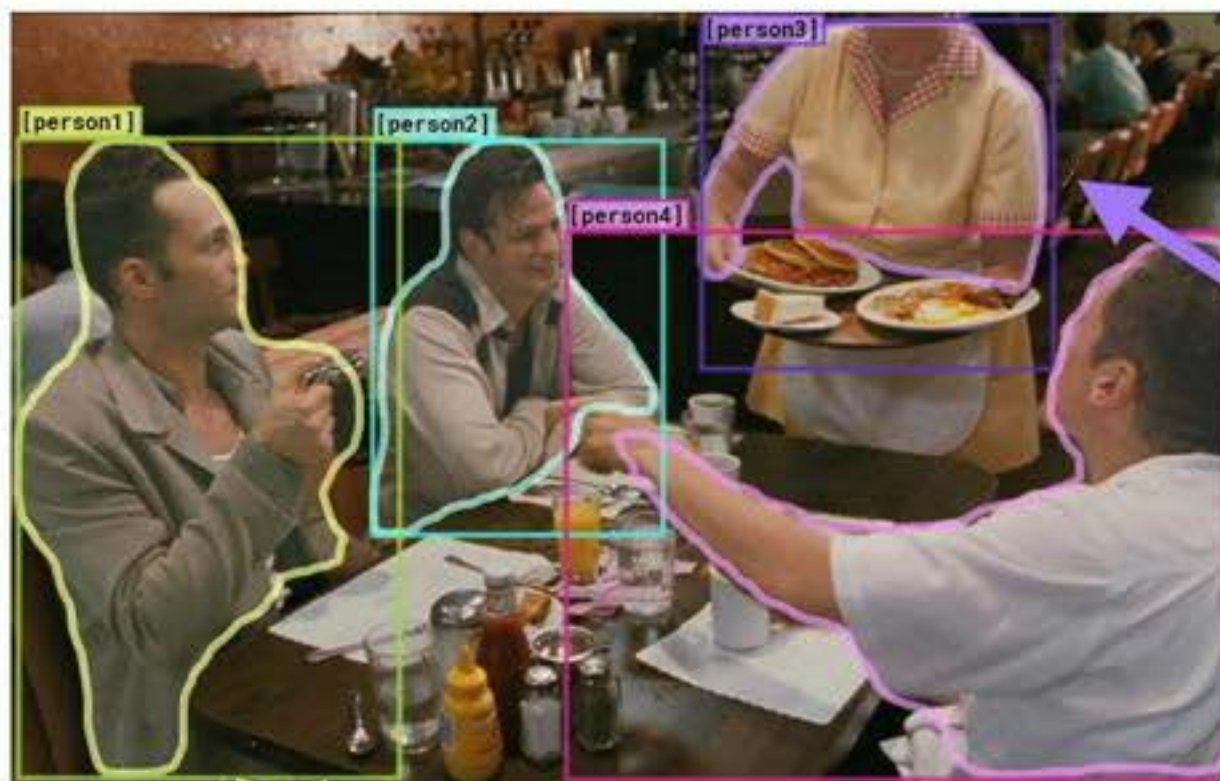
Objects

Part 1: Grounding

| Why | is | p4 | poin ting | ... |

Query

Why is [person4 🖼️] pointing at [person1 🖼️]?

| He | is | tell ing | p3 | ... |

Response Choice

He is telling [person3 🖼️] that [person1 🖼️] ordered pancakes.

Part 1: Grounding

Objects

**LSTM** → ← **LSTM** → ← **LSTM** → ← **LSTM** → ←

Why | is | p4 | poin ting | ...

*Query*

Why is [**person4** 🖼] pointing at [**person1** 🖼]?

**LSTM** → ← **LSTM** → ← **LSTM** → ← **LSTM** → ← ...

He | is | tell ing | p3 | ...

*Response Choice*

He is telling [**person3** 🖼]
that [**person1** 🖼]
ordered pancakes.

*Part 2: Contextualization*

Objects  ...

*Query*

Why is [**person4** ]
pointing at [**person1** ]?

*Response Choice*

He is telling [**person3** ]
that [**person1** ]
ordered pancakes.

*Part 2: Contextualization*

Objects

| | Why | is | | ... |
|---|---|---|---|---|
| He | | | | |
| is | | | | |
| telling | | | | |
| ... | | | | |

*Query*

**Why** is [**person4** ]

pointing at [**person1** ]?

*Response Choice*

He is telling [**person3** ]

that [**person1** ]

ordered pancakes.

*Part 2: Contextualization*

*Response Choice*

 is telling [person3 ]

that [person1 ]

ordered pancakes.

*Part 3: Reasoning*

| *Response* | He | is | telling | ... |
|---|---|---|---|---|
| *Attended Query* | [person4] | is | pointing | ... |
| *Attended Objects* | | | | ... |

*Response Choice*

is telling [person3 ]

that [person1 ]

ordered pancakes.

*Part 3: Reasoning*

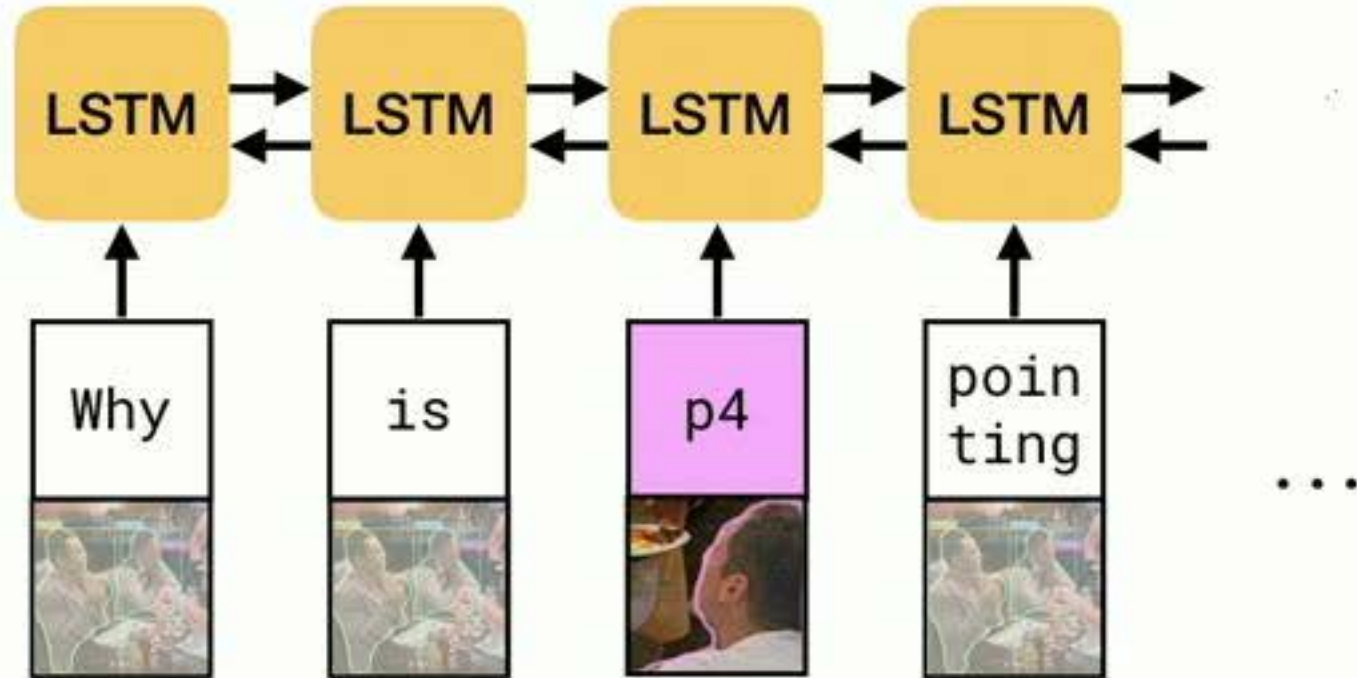**Max pool+multilayer perceptron**

| | LSTM$_1$ → | LSTM$_2$ → | LSTM$_3$ → |
|---|---|---|---|

| *Response* | He | is | telling | ... |
|---|---|---|---|---|
| *Attended Query* | [person4] | is | pointing | ... |
| *Attended Objects* | | | | ... |

# VCR Results

# VCR Results

**Accuracy**

■ Question -> Answer

| | Human | MLB | BERT | R2C |
|---|---|---|---|---|
| 100 | 91 | 46 | | |

Human    MLB    BERT    R2C

# VCR Results

# VCR Results

1. Why is **[person1]** pointing a gun at **[person2]** ?

a) **[person1]** wants to kill **[person2]** .    1.4%

b) **[person1]** and **[person3]** are robbing the bank and **[person2]** is the bank manager.    71.7%

c) **[person2]** has done something to upset **[person1]** .    18.7%

d) Because **[person2]** is **[person1]** 's daughter. **[person1]** wants to protect **[person2]** .    8.2%

1. Why is [person1] pointing a gun at [person2]?

a) [person1] wants to kill [person2].   1.4%

b) [person1] and [person3] are robbing the bank and [person2] is the bank manager.   71.7%

c) [person2] has done something to upset [person1].   18.7%

d) Because [person2] is [person1]'s daughter. [person1] wants to protect [person2].   8.2%

a) [person1] is chasing [person1] and [person3] because they just robbed a bank.   33.8%

b) Robbers will sometimes hold their gun in the air to get everyone's attention.   5.3%

c) The vault in the background is similar to a bank vault. [person3] is waiting by the vault for someone to open it.   49.1%

d) A room with barred windows and a counter usually resembles a bank.   11.7%

# visualcommonsense.com



**VCR — Visual Commonsense Reasoning**

Home | Paper | Download | Code | Explore | Leaderboard

## Visual Commonsense Reasoning

On VCR, a model must not only answer commonsense visual questions, but also provide a rationale that explains why the answer is true.

Read our paper on arXiv»

## Submitting to the leaderboard

Submission is easy! You just need to email Rowan with your predictions. Formatting instructions are below:

Submit to the leaderboard»

## What kinds of submissions are allowed?

The only constraint is that your system must predict the **answer first**, then the rationale. (The rationales were selected to be highly relevant to the correct Q,A pair, so they leak information about the correct answer.)

- To deter this, the submission format involves submitting predictions for each possible rationale, conditioned on each possible answer.
- A simple way of setting up the experiments (used in the paper) is to consider a task with *query* and four *response* choices. For Q->A the query is the question, and the response choices are the answers. For QA->R, the query is the question and answer, concatenated together, and the response choices are the rationales.

## Questions?

If it's not about something private, check out the google group below:

Get help via the google group»

## VCR Leaderboard

There are two different subtasks to VCR:

- **Question Answering** (Q->A): In this setup, a model is provided a question, and has to pick the best answer out of four choices. Only one of the four is correct.
- **Answer Justification** (QA->R): In this setup, a model is provided a question along with the correct answer, and it has to justify it by picking the best rationale out of four choices.

We combine the two parts with the **Q->AR** metric, in which a model only gets a question right if it answers correctly *and* picks the right rationale. Models are evaluated in terms of accuracy (%). How well will your model do?

| Rank | Model | Q->A | QA->R | Q->AR |
|---|---|---|---|---|
| | Human Performance<br>*University of Washington*<br>(Zellers et al. '18) | 91.0 | 93.0 | 85.0 |
| Feb 9, 2019 | CKRE<br>*Peking University* | **66.1** | **68.5** | **45.5** |
| 2<br>Nov 28, 2018 | Recognition to Cognition Networks<br>*University of Washington*<br>https://github.com/rowanz/r2c | 65.1 | 67.3 | 44.0 |
| 3<br>Nov 28, 2018 | BERT-Base<br>*Google AI Language*<br>*(experiment by Rowan)*<br>https://github.com/google-research/bert | 53.9 | 64.5 | 35.0 |
| 4<br>Nov 28, 2018 | MLB<br>*Seoul National University*<br>*(experiment by Rowan)*<br>https://github.com/jnhwkim | 46.2 | 36.8 | 17.2 |

visualcommonsense.com

VCR

VISUAL COMMONSENSE REASONING

Home | Paper | Download | Code | Explore | Leaderboard

## Visual Commonsense Reasoning

On VCR, a model must not only answer commonsense visual questions, but also provide a rationale that explains why the answer is true.

Read our paper on arXiv»

## Submitting to the leaderboard

Submission is easy! You just need to email Rowan with your predictions. Formatting instructions are below:

Submit to the leaderboard»

## What kinds of submissions are allowed?

The only constraint is that your system must predict the **answer first**, then the rationale. (The rationales were selected to be highly relevant to the correct Q,A pair, so they leak information about the correct answer.)

- To deter this, the submission format involves submitting predictions for each possible rationale, conditioned on each possible answer.
- A simple way of setting up the experiments (used in the paper) is to consider a task with *query* and four *response* choices. For Q->A the query is the question, and the response choices are the answers. For QA->R, the query is the question and answer, concatenated together, and the response choices are the rationales.

## Questions?

If it's not about something private, check out the google group below:

Get help via the google group»

## VCR Leaderboard

There are two different subtasks to VCR:

- **Question Answering** (Q->A): In this setup, a model is provided a question, and has to pick the best answer out of four choices. Only one of the four is correct.
- **Answer Justification** (QA->R): In this setup, a model is provided a question along with the correct answer, and it has to justify it by picking the best rationale out of four choices.

We combine the two parts with the **Q->AR** metric, in which a model only gets a question right if it answers correctly *and* picks the right rationale. Models are evaluated in terms of accuracy (%). How well will your model do?

| Rank | Model | Q->A | QA->R | Q->AR |
|---|---|---|---|---|
| | Human Performance<br>*University of Washington* | 91.0 | 93.0 | 85.0 |
| Feb 9, 2019 | CKRE<br>*Peking University* | 66.1 | 68.5 | 45.5 |
| 2<br>Nov 28, 2018 | Recognition to Cognition Networks<br>*University of Washington*<br>https://github.com/rowanz/r2c | 65.1 | 67.3 | 44.0 |
| 3<br>Nov 28, 2018 | BERT-Base<br>*Google AI Language*<br>*(experiment by Rowan)*<br>https://github.com/google-research/bert | 53.9 | 64.5 | 35.0 |
| 4<br>Nov 28, 2018 | MLB<br>*Seoul National University*<br>*(experiment by Rowan)*<br>https://github.com/jnhwkim | 46.2 | 36.8 | 17.2 |

# What's next?

# Future work

- VCR is a new testbed for commonsense visual reasoning!

- What kinds of new models will do well?

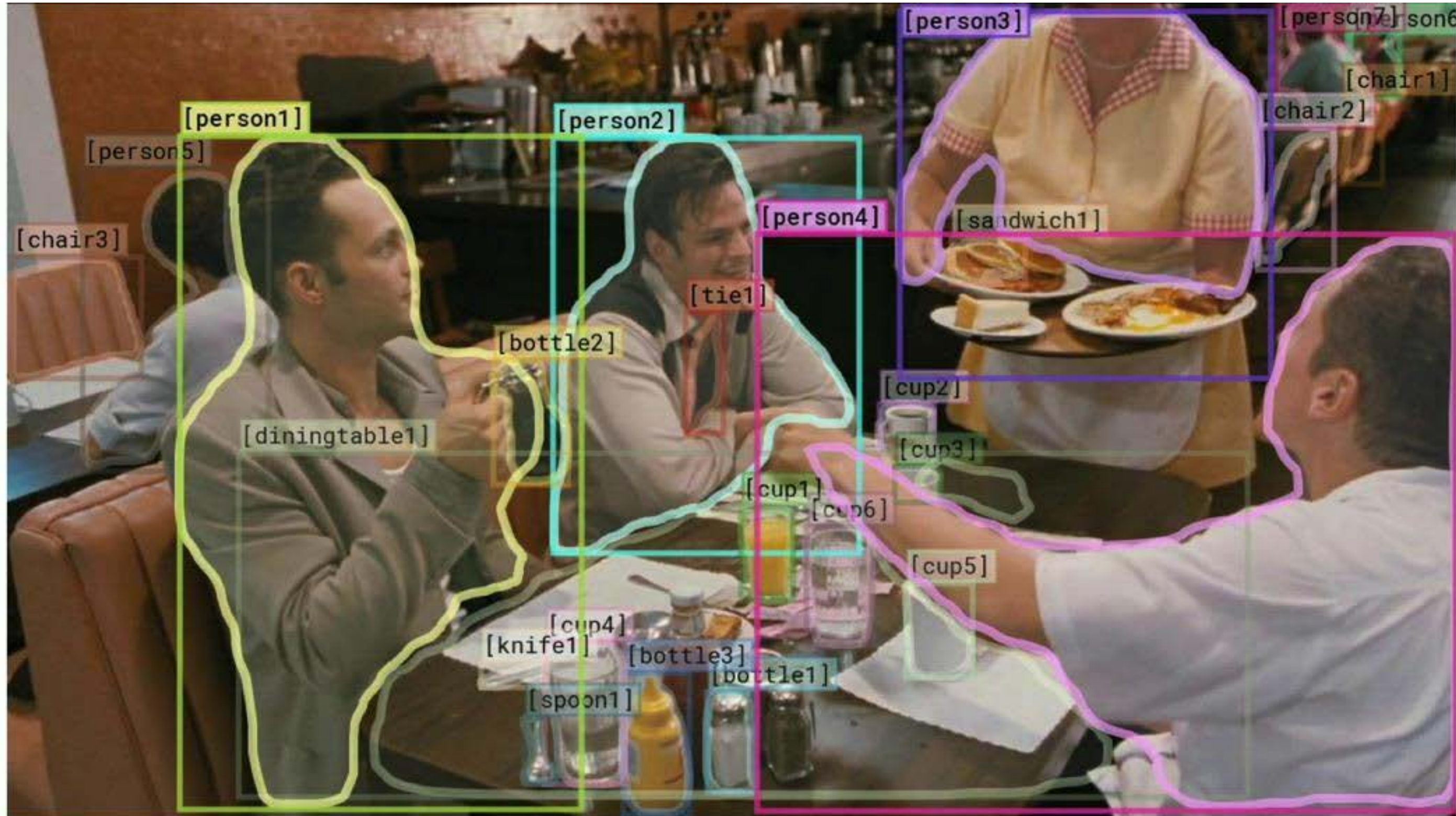# Summary

Thanks all!!
More at rowanzellers.com
twitter: @rown

Me     Yonatan Bisk     Roy Schwartz     Ali Farhadi     Yejin Choi

Why is [person4 🖼] pointing at [person1 🖼]?